

There are good reasons to randomize the errors of \vec{w}
 \Rightarrow the "randomized experiment" (1925, Fisher)

$\Rightarrow \vec{w} \sim \frac{1}{\binom{n}{n_T}} \mathbb{1}_{\sum w_i = n_T}$ is called the "completely randomized design".

$E[\vec{w}] = \frac{1}{2} \vec{1}_n$. This design has the added bonus that the bias disappears:

$$\begin{aligned} \beta_T &= \frac{2\vec{w}^T \vec{y} - \vec{1}^T \vec{y}}{\frac{n}{2}} = \frac{4}{n} \vec{w}^T (\beta_0 \vec{1}_n + \beta_u \vec{u} + \beta_T \vec{w} + \vec{\epsilon}) - \frac{2}{n} \vec{1}^T (\beta_0 \vec{1}_n + \beta_u \vec{u} + \beta_T \vec{w} + \vec{\epsilon}) \\ &= \frac{4}{n} \left(\beta_0 \underbrace{\vec{1}^T \vec{w}}_{\frac{n}{2}} + \beta_u \vec{u}^T \vec{w} + \beta_T \underbrace{\vec{w}^T \vec{w}}_{\frac{n}{2}} + \vec{w}^T \vec{\epsilon} \right) - \frac{2}{n} \left(\beta_0 \underbrace{\vec{1}^T \vec{1}_n}_n + \beta_u \vec{1}^T \vec{u} + \beta_T \underbrace{\vec{1}^T \vec{w}}_{\frac{n}{2}} + \vec{1}^T \vec{\epsilon} \right) \\ &= \cancel{2\beta_0} + \frac{4\beta_u}{n} \vec{u}^T \vec{w} + \cancel{2\beta_T} + \frac{4}{n} \vec{w}^T \vec{\epsilon} - \cancel{2\beta_0} - \frac{2\beta_u}{n} \vec{1}_n^T \vec{u} - \cancel{\beta_T} - \frac{2}{n} \vec{1}^T \vec{\epsilon} \end{aligned}$$

$$E_{\vec{\epsilon}}[\] = \beta_T + \frac{4\beta_u}{n} \vec{u}^T \vec{w} - \frac{2\beta_u}{n} \vec{1}_n^T \vec{u}$$

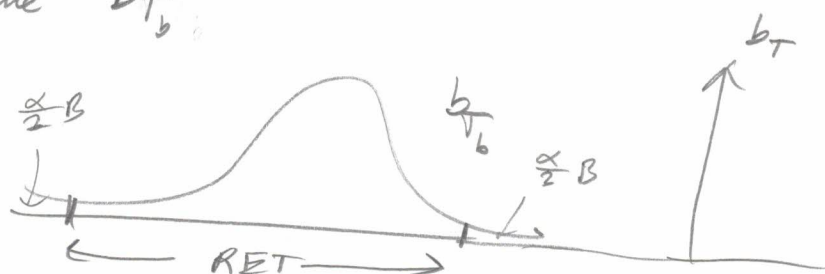
$$E_{\vec{w}}[\] = \beta_T + \frac{4\beta_u}{n} \vec{u}^T \left(\frac{1}{2} \vec{1}_n \right) - \frac{2\beta_u}{n} \vec{1}_n^T \vec{u} = \beta_T \checkmark$$

On average, over all ^{randomized} experiments, the PATE is measured unbiasedly

There is a whole other reason to employ randomization. Fisher called it a "random basis for inference". Here's how this works.

Let $H_0: Y_i(w_i=1) = Y_i(w_i=0) \forall i$ that is, the response value will be identical regardless of the manipulation

You run the the experiment with \vec{w}_{exp} and collect \vec{y} and compute b_T .
Then for $B = \text{large \#}$, you generate \vec{w} from \vec{W} and compute b_T .



Then you can draw null limits at the $\frac{\alpha}{2}$ and $1 - \frac{\alpha}{2}$ empirical quantiles.

If b_T from the actual experiment falls outside of RET \Rightarrow reject H_0 .

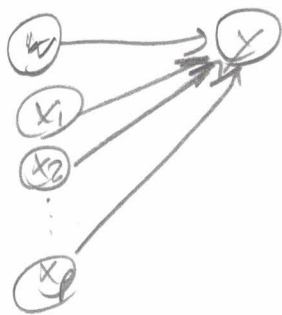
Note: no assumptions about the population model such as $\vec{E} \sim N_n(\vec{0}, \sigma^2 I_n)$ are necessary! This is a non-parametric test just like the perm. test.

\Rightarrow Randomization gives you:

- (1) insurance against unmeasured covariate or experimental imperfection
- (2) the ability to have a non-parametric test option

What about the least baseline covariates?

3



$$\text{let } \vec{Y} = \beta_0 \vec{1}_n + \beta_T \vec{w} + \beta_1 \vec{x}_1 + \dots + \beta_p \vec{x}_p + \vec{\epsilon}$$

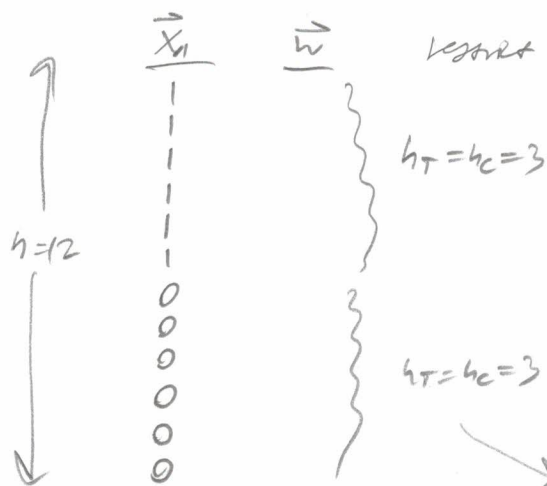
Similar to caption from before...

$$\Rightarrow B_T = \vec{Y}_T - \vec{Y}_C = \beta_T + \beta_1(\bar{x}_{1,T} - \bar{x}_{1,C}) + \beta_2(\bar{x}_{2,T} - \bar{x}_{2,C}) + \dots + \beta_p(\bar{x}_{p,T} - \bar{x}_{p,C}) + \vec{\epsilon}_T - \vec{\epsilon}_C$$

$$\Rightarrow E_{\vec{\epsilon}}[B_T] = \beta_T + \underbrace{\beta_1(\bar{x}_{1,T} - \bar{x}_{1,C}) + \dots + \beta_p(\bar{x}_{p,T} - \bar{x}_{p,C})}_{\text{Bias}}$$

We can reduce this bias by "restricting" the randomization to keep these terms small. This was noticed by Fisher immediately.

For example, let $X_1 = \text{gender}$ and let $n=12$



$$\Rightarrow \bar{x}_{1,T} = 0.5, \bar{x}_{1,C} = 0.5$$

$$\Rightarrow \beta_1(\bar{x}_{1,T} - \bar{x}_{1,C}) = 0 \text{ for all } \vec{w} !!$$

only choose \vec{w} 's that respect this constraint!

This is known as "blocking". The two blocks above are one for male, one for female.

If X_1 is continuous \Rightarrow binomial or cut in 3 or 4 even pieces (blocks)
then randomize with equal allocation within block

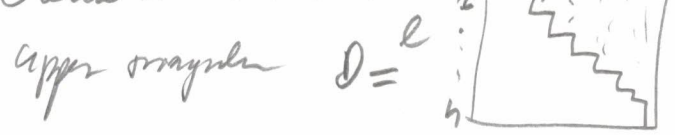
Blocking is the most common restricted design.

Horvath blocking doesn't work for p large as block size is $\frac{n}{2p} \rightarrow 0$ quickly. My research now shows that a good recommendation is

the pairwise ^{matcher} design: split the n subjects into $\frac{n}{2}$ pairs.

How? Define Euclidean distance: $d(\vec{x}_k; \vec{x}_l) = \sum_{j=1}^p (x_{kj} - x_{lj})^2$

Create distance matrix:



Then use the bipartite matching algorithm to create pairs that minimize the overall distance of total $\sum_{k,l} d_{k,l}$. Within pairs, randomize $w \in \{0,1\}$

You can think of this as $\frac{n}{2}$ blocks of size 2 each.

Another design that has gotten attention: rerandomization (Sturges' idea). Draw a \vec{w} . If $\sum_{j=1}^p (\bar{x}_{j,T} - \bar{x}_{j,C})^2$ is "large",

draw another \vec{w} until it's small i.e. $\leq D_{th}$, your hyperparameter.

For blocking, pairwise matching, rerandomization, you can use the rerandomization test as well... just make sure you draw \vec{w} 's that are legal within the restricted allocation set.

What if $y = \{C_1, C_2, \dots, C_K\}$ class from with K levels.

If $K=2$, lots simpler... assume $Y_i \stackrel{iid}{\sim} \text{Bern}(\phi(\vec{x}_i \vec{\beta}))$ where ϕ is a link function $\phi: \mathbb{R} \rightarrow (0,1)$.
 \Rightarrow Lik. Function:

$$L(\vec{\beta}; X, y) = \prod_{i=1}^n \phi(\vec{x}_i \vec{\beta})^{y_i} (1 - \phi(\vec{x}_i \vec{\beta}))^{1-y_i}$$
 if $\phi(u) = \frac{e^u}{1+e^u}$, logistic regression

If $K>2$, assume $Y_i \stackrel{iid}{\sim} \text{Multi}(1, \vec{\theta}_i) = \begin{pmatrix} 1 \\ y_{i1} y_{i2} \dots y_{iK} \end{pmatrix} \theta_1^{y_{i1}} \theta_2^{y_{i2}} \dots \theta_K^{y_{iK}}$
 let $\vec{\theta}_i = \begin{bmatrix} \phi(\vec{x}_i \vec{\beta}_1) \\ \phi(\vec{x}_i \vec{\beta}_2) \\ \vdots \\ \phi(\vec{x}_i \vec{\beta}_K) \end{bmatrix}$ and we $\vec{1}_K^T \vec{\theta}_i = 1$, $\phi(\vec{x}_i \vec{\beta}_K) = 1 - (\phi(\vec{x}_i \vec{\beta}_1) + \dots + \phi(\vec{x}_i \vec{\beta}_{K-1}))$

\Rightarrow Likelihood Function is:

$$L(\vec{\beta}_1, \vec{\beta}_2, \dots, \vec{\beta}_{K-1}; X, y) = \prod_{i=1}^n \phi(\vec{x}_i \vec{\beta}_1)^{y_{i1}} \phi(\vec{x}_i \vec{\beta}_2)^{y_{i2}} \dots \phi(\vec{x}_i \vec{\beta}_{K-1})^{y_{i,K-1}} (1 - (\phi(\vec{x}_i \vec{\beta}_1) + \dots + \phi(\vec{x}_i \vec{\beta}_{K-1})))^{y_{iK}}$$

If $\phi(u) = \frac{e^u}{1+e^u} \Rightarrow \text{multilogit model}$
 If $\phi(u) = \Phi(u) \Rightarrow \text{multinomial probit model}$ } most common

Inference? Use Wald test. What is the effect of x_1 ? There are $K-1$ effects:

$b_{k1}, b_{21}, \dots, b_{K-1,1}$ where each is a log-odds effect on prob of class k !

Adjusting for impact... as you get used to it.